

SocialBench: Sociality Evaluation of Role-Playing Conversational Agents

Hongzhan Chen¹, Hehong Chen², Ming Yan^{2*}, Wenshen Xu², Xing Gao²
Weizhou Shen¹, Xiaojun Quan^{1*}, Chenliang Li², Ji Zhang², Fei Huang², Jingren Zhou²

¹School of Computer Science and Engineering, Sun Yat-sen University, China

²Alibaba Group, China

¹chenhzh59@mail2.sysu.edu.cn, quanxj3@mail.sysu.edu.cn

²ym119608@alibaba-inc.com

Abstract

Large language models (LLMs) have advanced the development of various AI conversational agents, including role-playing conversational agents that mimic diverse characters and human behaviors. While prior research has predominantly focused on enhancing the conversational capability, role-specific knowledge, and stylistic attributes of these agents, there has been a noticeable gap in assessing their social intelligence. In this paper, we introduce SocialBench, the first benchmark designed to systematically evaluate the *sociality* of role-playing conversational agents at both individual and group levels of social interactions. The benchmark is constructed from a variety of sources and covers a wide range of 500 characters and over 6,000 question prompts and 30,800 multi-turn role-playing utterances. We conduct comprehensive evaluations on this benchmark using mainstream open-source and closed-source LLMs. We find that agents excelling in individual level does not imply their proficiency in group level. Moreover, the behavior of individuals may *drift* as a result of the influence exerted by other agents within the group. Experimental results on SocialBench confirm its significance as a testbed for assessing the social interaction of role-playing conversational agents. The benchmark is publicly accessible at <https://github.com/X-PLUG/SocialBench>.

1 Introduction

Recently, role-playing applications powered by LLMs, such as Character.AI¹, have gained significant attention. A growing number of research efforts have been dedicated to developing LLM-based role-playing conversational agents, aiming to mimic diverse characters and human behavior (Wang et al., 2023c; Shao et al., 2023; Tu et al., 2024; Zhou et al., 2023; Tian et al., 2023).

As an emerging and rapidly developing area, the evaluation of role-playing conversational agents is becoming increasingly important. Wang et al. (2023c) collected a role-specific instruction dataset and utilized Rouge-L and GPT 3.5 to assess the model’s role-specific knowledge and speaking style. Tu et al. (2024) proposed a Chinese benchmark and trained a reward model to measure the model’s conversational ability and character consistency and attractiveness. While these works mainly focus on evaluating the agent’s individual abilities to imitate the character’s role-specific knowledge or speaking style, this study aims to explore and measure the *social interaction* of role-playing conversational agents, another pivotal dimension for assessing how role-playing agents perceive and behave in a social interaction environment.

Therefore, we introduce SocialBench, the first evaluation benchmark designed to systematically assess the social interaction of role-playing conversational agents. As introduced in (Troitzsch, 1996), the agent society represents a complex system comprising individual and group social activities. Following this definition, SocialBench assesses the sociality metrics at both the individual and group levels, as shown in Figure 1. At the individual level, the agent should possess the basic social intelligence as individuals, such as self-awareness on role description (Wang et al., 2023c; Tu et al., 2024; Shen et al., 2023), emotional perception on environment (Hsu et al., 2018), and long-term conversation memory (Zhong et al., 2023). Each of these aspects contributes to the nuanced understanding of how the agents manifest their individual social behaviors. Moreover, we further examine the social intelligence of the role-playing agents within group social interactions, which require the agents to possess certain social preferences towards group dynamics (Leng et al., 2023).

SocialBench is carefully constructed from diverse English and Chinese books, movies, and

* Corresponding authors.

¹<https://beta.character.ai>

novels, covering a wide range of 500 characters and 6,000 questions, and 30,800 multi-turn role-playing utterances. Specifically, we design a three-step construction pipeline for SocialBench. Firstly, we collect diverse role profiles from common web sources. Secondly, GPT-4 is employed to extract dialogue scenes, individual and group-level social conversations, as well as multi-choice questions. Thirdly, we conduct a series of pre-processing and manual labeling to ensure the quality of the benchmark. We conduct comprehensive evaluations on SocialBench using mainstream open-source and closed-source LLMs to inspire future research.

2 Related Work

2.1 Role-Playing Agents

Leveraging the powerful capabilities of open-source foundational models, numerous efforts have emerged to develop models specifically tailored for role-playing tasks. These approaches can be categorized based on training paradigms: 1) Supervised fine-tuning (SFT). Li et al. (2023); Wang et al. (2023c); Tu et al. (2023) involved constructing specialized persona training corpus while performing fine-tuning on it to enable the agents to acquire capabilities of role-playing. 2) Integration of offline reinforcement learning. Shea and Yu (2023) combined role-playing model training with importance sampling strategies. 3) Incorporation of retrieval-enhanced methods. Salemi et al. (2023) combined role-playing model training with retrieved information to enhance the capabilities of agents in role-playing. (Shao et al., 2023) introduced a experience upload method, to test the model’s effectiveness on memorizing the character knowledge, values and personality.

2.2 Role-Playing Benchmarks

With the rapid development of role-playing agents, there has been a corresponding increase in evaluation datasets. Current evaluation datasets mostly focus on the alignment of role-playing agents with regards to role style and role knowledge. In terms of role style, Tu et al. (2024) and Wang et al. (2023c) investigate whether models can generate responses consistent with the style of the given role. Agents need to grasp different speaking styles for different roles. Regarding role knowledge, Shen et al. (2023) particularly focuses on the role knowledge of role-playing models, including the characters’ experiences and social relationships. Tu et al. (2024) and

Wang et al. (2023c) also address aspects of role knowledge, such as role knowledge illusions. Additionally, Wang et al. (2023a) and Tu et al. (2024) introduce psychological theories like the Big Five and MBTI to evaluate role-playing agents. While previous work mainly focuses on testing the abilities of agents on imitating the character’s role-specific knowledge or speaking style, SocialBench introduces the first-ever evaluation benchmark for the sociality of role-playing conversational agents encompassing both individual and group level.

2.3 Agent Society

Previous benchmarks have primarily focused on single-agent scenarios, leaving the more complex multi-agent scenarios underexplored. Similar to humans, agents are capable of engaging in intricate social interactions, resulting in the formation of an agent society (da Rocha Costa, 2019). Recently, LLM-based agents demonstrate complex social behaviors, where cooperation and competition coexist (Xu et al., 2023). These sophisticated behaviors intertwine to shape social interactions (Gao et al., 2023). SocialBench follows the framework defined by Nigel Gilbert and Troitzsch (1997); Leng et al. (2023), where behaviors in agent societies are divided into individual and group-level activities, to study the social intelligence of role-playing agents within social interactions.

3 Sociality of Role-Playing Agent

The role-playing agent is designed to engage in conversations with users by imitating predefined characters. Given the character profile and social context, the sociality of role-playing agents focuses on imitating typical human social interactions from individual level to group level.

3.1 Individual Level

At the individual level, the role-playing social agents manifest through various capabilities, which collectively contribute to their ability to interact within a social context. These capabilities form the foundation of the agent’s social behavior.

Self-Awareness on Role Description involves understanding not only the role’s knowledge (Shen et al., 2023), but also the role’s distinct behavioral style (Zhou et al., 2023; Wang et al., 2023b). This self-awareness enables the agent to maintain consistency with its designated role.

Emotional Perception on Environment enables agents to acquire high-level feeling percep-

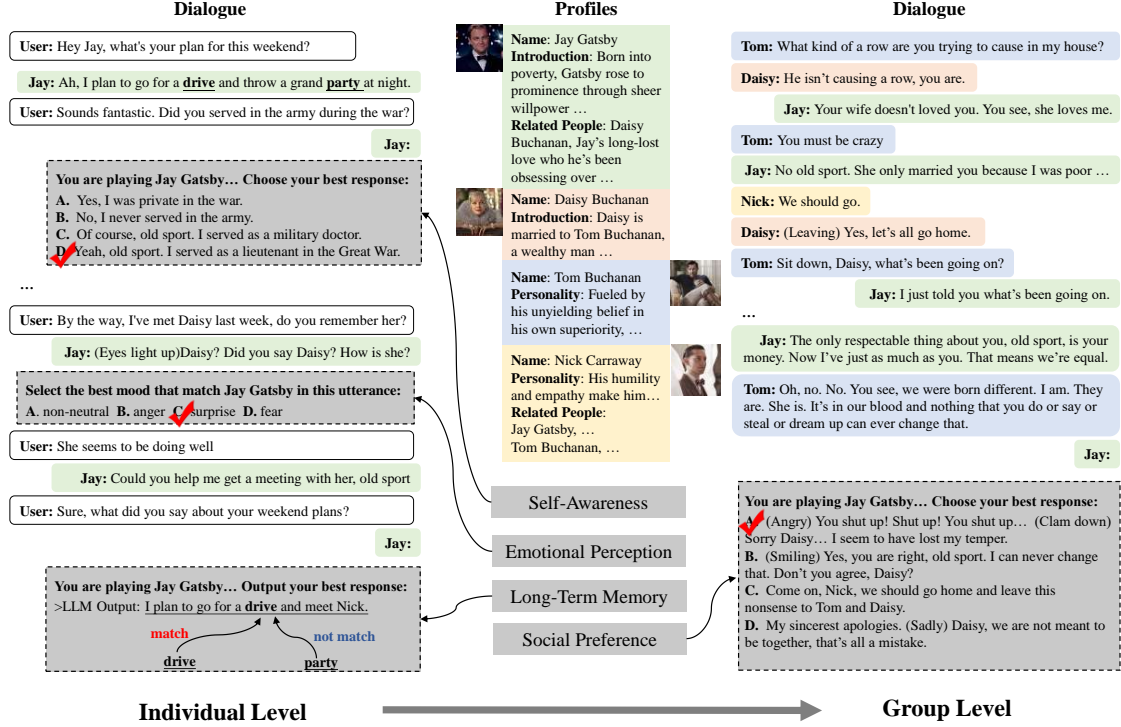


Figure 1: An example from SocialBench, which is partially constructed from the film “The Great Gatsby”.

tion for effective social interactions (Hsu et al., 2018). Agents endowed with sophisticated emotional intelligence, such as situation understanding and emotion detection, can perceive and respond to the emotions of others, facilitating smoother communication and relationship-building.

Long-Term Conversation Memory is crucial for conversational agents (Shao et al., 2023; Zhong et al., 2023). By memorizing previous dialogue content and aligning with their statements accordingly, role-playing agents demonstrate reliability, enhancing the quality of their social engagements.

3.2 Group Level

Individuals within group conversation may be influenced by the group member interactions, thus demonstrate more sophisticated social behaviors towards group dynamics. It represents a higher calling for the sociality of role-playing agent.

Social Preference towards Group Dynamics. As a group member, it is natural to navigate diverse group conversation scenarios: acting as a leader to control the pace of conversation, serving as a mediator when conflicts arise among the group, or considering others’ perspectives during discussion, which shows its internal social preference (Leng et al., 2023) towards group dynamics. Furthermore,

within society, not all behaviors are inherently positive for the group, and some may be neutral or even negative (Xi et al., 2023). Social agents need to exhibit and keep their pre-designed social preference or group identity when confronted with diverse and more sophisticated group conversations.

4 SocialBench

In this section, we introduce the construction process of SocialBench, as illustrated in Figure 2.

4.1 Profile Collection

A role profile defines the character style, knowledge, emotions, and social preference of a role-playing agent. We gather profiles for role-playing agents from various sources including novels, scripts, online platforms such as CharacterAI² and Fandom³, and automatic generation via GPT-4-Turbo prompting. To ensure diversity, we construct profiles based on various character types and personality traits by combining the existing categorizations in online platforms and research work (Shen et al., 2023; Gunkel, 1998). Figure 3 illustrates the distribution of personality traits for roles within SocialBench. It demonstrates our approach

²<https://beta.character.ai>

³<https://www.fandom.com>

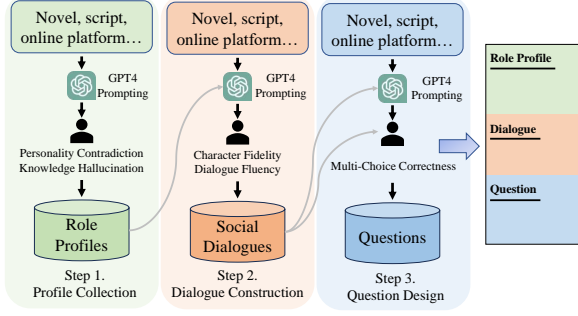


Figure 2: The three-step dataset construction pipeline of SocialBench.

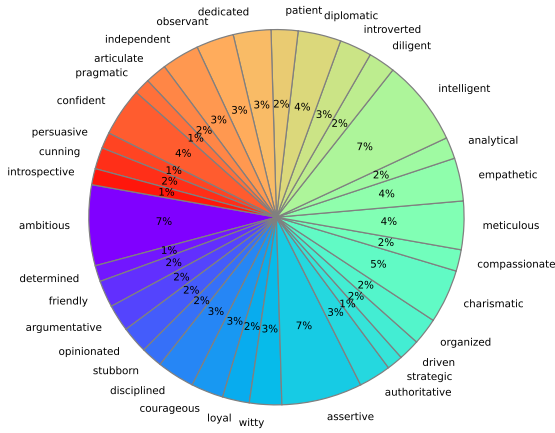


Figure 3: Personality traits distribution in SocialBench.

of ensuring both category diversity and balanced proportions across all categories. The details regarding profile collection in SocialBench can be found in Appendix D.1.

4.2 Dialogue Construction

The dialogue construction adheres to two principles: *dialogue fluency*, which ensures natural and coherent conversations; and *character fidelity*, meaning all characters in the dialogue must adhere to their respective personas. We employ four dialogue construction methods: 1) Extracting from novels and scripts: We gather novels and scripts and extract high-quality dialogue data. 2) Collecting from online role-playing platforms: We collect authentic user dialogue data from online role-playing platforms. 3) Conducting role-playing tasks between users and general LLMs: We prompt general LLMs like GPT-4-Turbo to role-play characters and engage users to generate dialogue data. 4) Fully automatic self-dialogue generation with

general LLMs: We task general LLMs like GPT-4-Turbo to role-play and engage in self-dialogue for data collection. Prompts for extracting dialogues can be found in Appendix B.1.

4.3 Question Design

Based on the constructed dialogues, we employ different methods for designing questions tailored to different dimensions within SocialBench.

For Self-Awareness: This includes two subcategories: self-awareness on role style (SA Style) and self-awareness on role knowledge (SA Know.). Utterances from the original dialogue are selected as correct answers because they have been manually verified to conform to the corresponding role style and role knowledge. For SA Style, we choose styles contradicting the character as negative options, such as rephrasing the original sentence using a different tone. For SA Know., we alter correct answers to be inconsistent with the facts mentioned in the original sentence (e.g., time, location) as negative options.

For Emotional Perception: We construct questions related to situational understanding (EP Situ.) and emotion detection (EP Emo.) based on professional exam questions and relevant open-source datasets (Chen et al., 2022; Hsu et al., 2018; Garbowicz, 2021). For EP Situ., we design questions that require agents to analyze the psychological state of the speaker and identify the causes of this state. In this dimension, we design multiple-choice questions with multiple correct answers. For EP Emo., we design questions that require agents to analyze the current speaker’s emotions, such as happiness or sadness, with the correct emotion serving as the correct answer and incorrect emotions as negative options. We further use expert annotations or existing labels to create correct answers, while negative options are constructed through manual collection and generation by GPT-4-Turbo.

For Conversation Memory: This category includes two subcategories: short-term conversation memory (CM Short) and long-term conversation memory (CM Long). In a dialogue, the user initially poses a question to the agent. After several rounds of conversation, the user repeats the same question, expecting the agent to provide a consistent response. We employ keyword matching, where the agent is required to include the previously mentioned keywords in their subsequent response. For CM Short, we prompt the agent to recall keywords discussed within 40 utterances, while

for CM Long, we prompt the agent to recall keywords discussed over 40 utterances. We evaluate how many of these keywords are recalled.

For Social Preference: We design questions for three social behavior preferences: positive (Pos.), neutral (Neu.), and negative (Neg.). Group dialogues typically consist of social interactions involving 2 to 10 characters. We analyze the social preference of an agent and identify behaviors aligning with its preference in the dialogues as correct answers. For instance, behaviors like cooperation and coordination are deemed consistent with the preferences of an agent inclined towards positive social interactions, and thus, are designated as correct answers. Behaviors contradicting its social preference serve as negative options, such as behaviors reflecting a negative social preference, including refusal to cooperate or engage in competition. Other agents in group dialogues also have their own social behavior preferences, which are reflected in their profiles or demonstrated through their social interactions.

4.4 Dataset Validation

The validation stage includes two parts: dataset pre-validation and post-validation. Throughout this process, we undergo multiple iterations of rigorous manual screening, annotation, and refinement.

4.4.1 Dataset Pre-Validation

Profile Verification: After profile collection, we assess personality contradictions and knowledge hallucinations in profiles to ensure character accuracy. We manually review and modify any erroneous descriptions in profiles, while also ensuring the exclusion of specific personal information such as phone numbers and home addresses.

Dialogue Verification: Our focus is on ensuring dialogues adhere to principles of *dialogue fluency* and *character fidelity*. For fluency, we manually inspect dialogues for contextual coherence and natural expression. For fidelity, we analyze the speaker’s profile to verify if the utterance aligns with the character’s speaking style and behavior. Dialogues that do not meet requirements undergo manual correction.

Question Verification: For multiple-choice questions, we invite three different annotators to label each question. If all three annotators deem the question valid and agree on the answer, it is considered valid. For open-domain generation questions, we verify the correctness and validity of keywords

provided. Invalid questions are either modified by experts or discarded.

4.4.2 Dataset Post-Validation

We undergo the post-validation process after completing each round of dataset. Different dimensions require different validation strategies.

Validation for Self-Awareness: We focus on examining knowledge-related errors in the questions and options, particularly those generated by LLMs that may give rise to knowledge hallucinations. We remove questions that do not meet the requirements, while options that do not meet the requirements will be flagged for correction in the subsequent iteration.

Validation for Emotional Perception: Some of the questions we collect are sourced from professional psychology exams, which may include highly specialized content not conducive to assessing the basic abilities of role-playing agents. Therefore, we filter out samples that are too focused on psychology-specific knowledge, retaining those that are more general and fundamental for role-playing agents.

Validation for Conversation Memory: In this dimension, we’ve observed that questions containing pronouns (such as "him," "it," "she") often result in unclear or ambiguous references to preceding context. Therefore, we remove questions containing pronouns to prevent ambiguity. Additionally, we assess the validity of extracted keywords to ensure they are proper nouns, thereby avoiding mismatches caused by different verb tenses.

Validation for Social Preference: We find that the options within this dimension may exhibit similarities, making it difficult to distinguish the correct option from the negative ones. To reduce difficulty, we manually examine the similarity between options for each sample. For options with excessively high similarity, we increase the differentiation between negative options and the correct answer. For instance, if the correct option has a positive social preference, we select negative social preference content with significantly different characteristics as negative options.

5 Experiment Settings

In this section, we show the statistic of SocialBench. Then we introduce the metrics along with the evaluation LLMs.

Metrics	Individual Level						Group Level		
	SA Style	SA Know.	EP Situ.	EP Emo.	CM Short	CM Long	Pos.	Neu.	Neg.
	Acc_{single}	Acc_{single}	$Acc_{multiple}$	Acc_{single}	$Cover$	$Cover$	Acc_{single}	Acc_{single}	Acc_{single}
#Questions	1,063	1,408	193	1,016	773	1,348	586	724	606
Avg Utterances	17.9	9.4	1.0	6.4	23.9	76.7	15.6	16.1	16.0
Avg Tokens per Utterance	32.6	66.7	286.3	23.0	37.6	41.2	38.8	38.7	42.0
Avg Characters per Question	2	2	N/A	N/A	2	2	6.3	6.5	6.7

Table 1: Metrics and statistics of SocialBench. There are a total of 500 roles, comprising 6,000 questions and 30,800 utterances in SocialBench.

5.1 Dataset Statistic

We show the statistic of SocialBench in Table 1 and the distribution of dialogue tokens length in Figure 4. SocialBench encompasses individual level and group level. There are six subcategories in individual level: self-awareness on role style (SA Style), self-awareness on role knowledge (SA Know), situational understanding (EP Situ.), emotion detection (EP Emo.), short-term conversation memory (CM Short), and long-term conversation memory (CM Long). Group level consists of three social preference categories: positive (Pos.), neutral (Neu.), and negative (Neg.). There are a total of 500 roles, comprising 6,000 questions and 30,800 utterances in SocialBench.

5.2 Evaluation Metrics

Most of the previous methods (Wang et al., 2023c; Shao et al., 2023) for role-playing applications rely on GPT-3.5 or GPT-4 for evaluation, which may suffer from questionable accuracy on the role-playing scenario and costly API usage. We follow the popular benchmark MMLU (Hendrycks et al., 2020) and C-Eval (Huang et al., 2023), and prompt for automatic and fast evaluation free from LLMs. SocialBench utilizes fully automatic evaluation metrics, employing both multiple-choice and open-domain generation questions.

For single-answer questions, we calculate the accuracy (Acc_{single}) using the following formula:

$$Acc_{single} = \frac{\text{Number of correctly chosen options}}{\text{Total number of single-answer questions}} \quad (1)$$

For multiple-answer questions, we calculate the accuracy ($Acc_{multiple}$) using the following formula:

$$Acc_{multiple} = \sum_i^N \frac{\text{Score}_i}{\text{MaxScore}_i}, \quad (2)$$

where N is the total number of multiple-answer questions. Score_i is the score obtained for the i th question, considering both correct and partially correct options chosen. MaxScore_i is the maximum

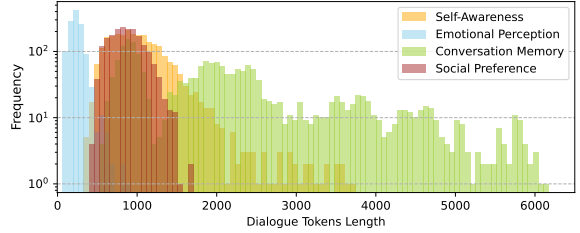


Figure 4: Distribution of dialogue tokens across four dimensions in SocialBench, based on tokenizer of Qwen.

achievable score for the i th question. For detailed information on metrics related to multiple-answer questions, please refer to Appendix C.

For open-domain question, we calculate the keyword coverage rate ($Cover$). SocialBench provides a keyword set $\mathbf{A}_{keywords} = \{k_1, k_2, \dots, k_n\}$. Given the keywords set mentioned in the response $\mathbf{R}_{keywords}$, we compute:

$$Cover(\mathbf{R}) = \frac{\text{len}(\mathbf{A}_{keywords} \cap \mathbf{R}_{keywords})}{\text{len}(\mathbf{A}_{keywords})}, \quad (3)$$

where $Cover(\cdot)$ quantifies the proportion of keywords mentioned in the response \mathbf{R} relative to the keywords identified in the \mathbf{A} .

The metrics utilized for different dimensions in SocialBench are listed in Table 1.

5.3 Models

We conduct evaluation on the current mainstream open-source and closed-source LLMs. For evaluation of open-source LLMs, we select chat version of LLaMA-2-7B/13B/70B (Touvron et al., 2023), instruction version of Mistral-7B (Instruct-V0.2) (Jiang et al., 2023), and chat versions of Qwen-7B/14B/72B (Bai et al., 2023). For evaluation of closed-source LLMs, we choose Minimax (abab5.5s-chat and abab6-chat)⁴, GLM (CharGLM-3 and GLM-3-Turbo) (Zhou et al., 2023), Baichuan (Baichuan-NPC-Turbo and

⁴<https://api.minimax.chat/>

Models (Max Length)	Individual Level						Group Level			Avg
	SA Style	SA Know.	EP Situ.	EP Emo.	CM Short	CM Long	Pos.	Neu.	Neg.	
Open-Source Models										
LLaMA-2-7B-Chat (4k)	48.76	51.23	31.23	28.91	25.38	21.89	44.98	24.19	27.67	33.80
LLaMA-2-13B-Chat (4k)	57.62	65.51	37.12	32.56	30.43	29.82	66.38	42.25	26.27	43.11
LLaMA-2-70B-Chat (4k)	67.61	70.78	35.74	38.47	45.57	26.74	69.87	45.29	39.37	48.83
Mistral-7B (8k)	50.12	61.17	36.48	31.72	31.78	25.42	65.67	46.34	28.96	41.96
Qwen-7B-Chat (32k)	66.44	71.16	41.68	40.68	67.45	53.45	75.61	52.78	43.11	56.93
Qwen-14B-Chat (32k)	77.06	86.15	45.71	43.78	65.32	51.37	78.32	58.25	59.21	62.80
Qwen-72B-Chat (32k)	83.87	90.64	53.10	52.89	<u>83.29</u>	73.15	<u>91.53</u>	73.44	63.82	73.97
Closed-Source Models										
GPT-4-Turbo (128k)	<u>84.57</u>	<u>93.11</u>	<u>56.48</u>	<u>53.05</u>	81.39	80.11	89.73	<u>81.69</u>	<u>75.10</u>	<u>77.25</u>
GPT-3.5-Turbo (16k)	73.17	73.82	52.44	45.49	73.03	59.72	81.59	76.79	54.16	65.58
Qwen-Max (8k)	82.04	93.34	61.14	52.36	76.45	72.65	87.22	72.14	52.19	72.17
Xingchen-Plus (8k)	85.43	91.6	55.44	60.73	82.43	80.69	94.27	86.69	77.26	79.39
Baichuan-NPC-Turbo (unknown)	53.69	61.67	52.14	43.34	76.47	22.40	62.09	48.91	34.59	50.59
Baichuan-2-Turbo (unknown)	77.75	83.35	55.7	47.38	80.11	78.91	87.37	74.71	68.50	72.64
CharGLM-3 (unknown)	74.70	79.41	26.23	41.27	81.16	68.29	84.40	70.45	36.36	62.47
GLM-3-Turbo (128k)	77.85	84.62	35.58	<u>53.05</u>	74.64	71.68	84.41	67.47	54.55	67.09
Minimax-abab5.5s-chat (8k)	36.09	42.11	28.15	47.97	29.55	19.30	44.59	41.04	22.45	34.58
Minimax-abab6-chat (32k)	82.92	87.45	35.90	51.38	83.60	<u>80.26</u>	89.12	79.55	74.65	73.87

Table 2: Main results from SocialBench. Best performances are shown in **bold**, while suboptimal ones underlined.

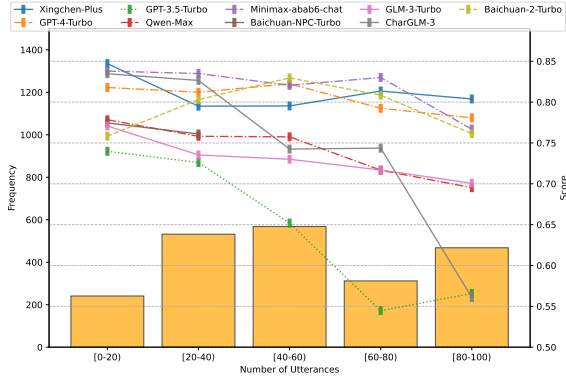


Figure 5: Performance w.r.t the number of utterances.

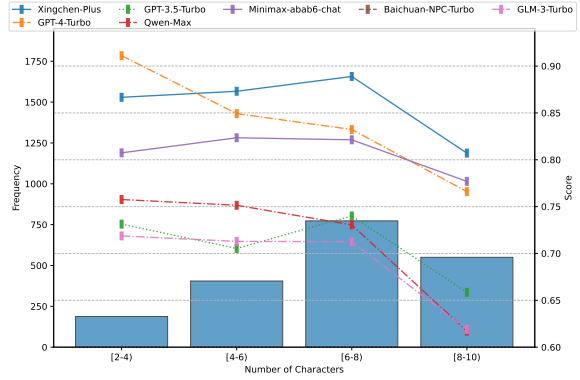


Figure 6: Performance w.r.t number of group members.

Baichuan-2-Turbo)⁵, Qwen-Max⁶, GPT-4-Turbo (OpenAI, 2023), GPT-3.5-Turbo (OpenAI, 2022), and Xingchen-Plus⁷.

6 Results and Analysis

In this section, we evaluate mainstream open-source and closed-source LLMs, while also analyzing the experimental results.

6.1 Overall Results

As presented in Table 2, the performance of closed-source models tends to surpass open-source models. Moreover, models specifically designed for role-playing, such as Xingchen-Plus, outperform others. While the general model GPT-4-Turbo also demonstrates impressive performance. However,

role-playing agents, like Baichuan-NPC-Turbo, CharGLM-3 and Minimax-abab5.5s, tend to underperform compared to their general counterparts, such as Baichuan-2-Turbo, GLM-3-Turbo and Minimax-abab6-chat. We find that they are biased towards character-based dialogues, leading to poorer understanding and compliance with instructions. It’s essential for role-playing agents to maintain character-based dialogue abilities and general instruction-following capabilities. At the individual level, dimensions such as SA Style, SA Know., and CM Short are well-performed by most models. However, some models tend to exhibit poor performance in EP Situ., EP Emo., and CM Long. At the group level, most models perform poorly due to the complexity of group dynamics. While models generally align well with tendencies towards positive social preference, there is a notable absence of necessary abilities to embody neutral and negative social preferences, which are also important for role-playing agents.

⁵<https://npc.baichuan-ai.com/index>

⁶<https://help.aliyun.com/zh/dashscope/developer-reference/api-details>

⁷<https://xingchen.aliyun.com/>

6.2 Conversation Memory for Role-Playing

Conversation memory capability is crucial for role-playing agents. We investigate the memory capacity of role-playing agents across different conversation lengths, measured by the number of utterances in the dialogue. We analyze the distribution of utterance counts in the conversation memory dimension of SocialBench. As illustrated in Figure 5, there is a declining trend in memory capability for some models, such as GPT-3.5-Turbo and CharGLM-3, as conversation length increases. When the number of utterances in the dialogue exceeds 80 rounds, most role-playing agents exhibit a noticeable performance decline. This finding showcases the limitations of current role-playing agents in handling extremely long-term memory and highlights potential areas for improvement.

6.3 Impact of Group Dynamics Complexity

We measure complexity of group dynamics by the number of group members, where a greater number denotes more intricate group dynamics. We analyze the distribution of the number of participating roles in group-level questions. As shown in Figure 6, with increasing complexity of group dynamics, the performance of all role-playing agents shows a downward trend. This can be interpreted as the interactions among a greater number of participants forming more complex group dynamics. We find that excelling in simple group dynamics does not necessarily imply their proficiency in more complex group dynamics. For example, models like GLM-3-Turbo and GPT-4-Turbo perform well in simple group dynamics, but this doesn’t guarantee strong performance in complex group dynamics. However, models like Xingchen-Plus and Minimax-abab6-chat, which are specially designed and trained with multi-turn role-playing data, could also demonstrate proficiency in handling complex group dynamics.

6.4 Impact of Group Dynamics Polarity

It is important for role-playing agents to maintain designed social preferences under the influence of varying group dynamics. The group dynamics polarity is defined as the majority social preference of group members. For instance, positive group dynamics imply that the majority of members exhibit positive social preference. For an individual with a specific social preference, different polarities of group dynamics may have various impacts. We

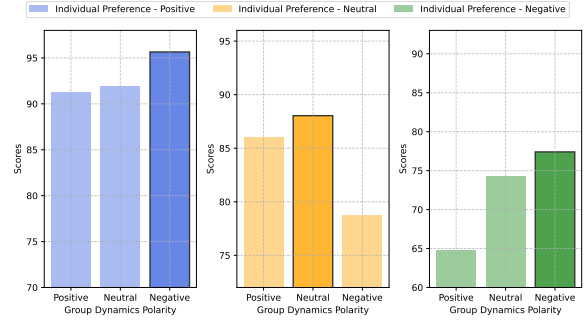


Figure 7: Performance of Xingchen-Plus under different group dynamics polarities on a subset of group data.

study the performance of individuals under different polarities of group dynamics, by analyzing a subset of group data in SocialBench. As shown in Figure 7, we find that individuals with neutral and negative social preferences perform optimally within their corresponding group polarities (i.e., neutral and negative group polarities). However, they are susceptible to the influence of group dynamics with different polarities and undergo a phenomenon termed as *preference drift*, leading to deviation from their original designed behaviors, as indicated by the decline of performance. Nevertheless, individuals with positive social preference appear to be more resilient to the preference drift, performing better across all group polarities. Especially, they excel in negative group polarity. This phenomenon can be termed as *social facilitation* (Guerin, 2010) in sociology. We hypothesize that negative group further motivates individuals to engage in behaviors advantageous to the group.

7 Conclusion

In this paper, we introduce SocialBench, the first evaluation benchmark designed to systematically assess the social intelligence of role-playing conversational agents at both individual and group levels. We construct diverse question prompts on a wide range of characters covering comprehensive dimensions, including self-awareness on role description, emotional perception on environment, long-term conversation memory, and social preference towards group dynamics. Moreover, rigorous human verification ensure questions’ difficulty and validity. We evaluate over 10 mainstream LLMs on SocialBench and provide in-depth analysis. While role-playing agents demonstrate satisfactory performance at the individual level, we find that their social interaction capabilities at the group level remain deficient. We hope this finding may inspire future research in this field.

Limitations

While SocialBench provides a comprehensive evaluation framework for assessing the sociality of role-playing conversation agents, there are several limitations to consider. 1) Social interactions, particularly within group settings, are inherently complex and nuanced. Despite our efforts, further research is needed to fully understand and capture the intricacies of these interactions. 2) The number of role-playing agents in group scenarios is relatively limited in our benchmark. Increasing the diversity and quantity of agents would provide a more comprehensive evaluation of the agents' social abilities and dynamics within groups. 3) Our dataset may contain some biased content, posing a risk of improper use. These limitations highlight areas for future research and development in the evaluation of social intelligence in role-playing agents.

References

- Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenhang Ge, Yu Han, et al. 2023. [Qwen technical report](#). *ArXiv*, abs/2309.16609.
- Yirong Chen, Weiquan Fan, Xiaofen Xing, Jianxin Pang, Minlie Huang, Wenjing Han, Qianfeng Tie, and Xiangmin Xu. 2022. [Cped: A large-scale chinese personalized and emotional dialogue dataset for conversational ai](#). *ArXiv*, abs/2205.14727.
- Antônio Carlos da Rocha Costa. 2019. *A Variational Basis for the Regulation and Structuration Mechanisms of Agent Societies*.
- Chen Gao, Xiaochong Lan, Zhi jie Lu, Jinzhu Mao, Jing Piao, Huandong Wang, Depeng Jin, and Yong Li. 2023. [S3: Social-network simulation system with large language model-empowered agents](#). *ArXiv*, abs/2307.14984.
- Krzysztof Garbowicz. 2021. [Dilbert2: Humor detection and sentiment analysis of comic texts using fine-tuned bert models](#).
- Xiaochang Gong, Qin Zhao, Jun Zhang, Ruibin Mao, and Ruifeng Xu. 2020. [The design and construction of a Chinese sarcasm dataset](#). In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 5034–5039, Marseille, France. European Language Resources Association.
- Bernard Guerin. 2010. Social facilitation. *The Corsini encyclopedia of psychology*, pages 1–2.
- Patrick Gunkel. 1998. [Human kaleidoscope](#).
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2020. Measuring massive multitask language understanding. *arXiv preprint arXiv:2009.03300*.
- Chao-Chun Hsu, Sheng-Yeh Chen, Chuan-Chun Kuo, Ting-Hao Huang, and Lun-Wei Ku. 2018. [Emotion-Lines: An emotion corpus of multi-party conversations](#). In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*, Miyazaki, Japan. European Language Resources Association (ELRA).
- Yuzhen Huang, Yuzhuo Bai, Zhihao Zhu, Junlei Zhang, Jinghan Zhang, Tangjun Su, Junteng Liu, Chuancheng Lv, Yikai Zhang, Jiayi Lei, et al. 2023. [C-eval: A multi-level multi-discipline chinese evaluation suite for foundation models](#). *arXiv preprint arXiv:2305.08322*.
- Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, L  lio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timoth  e Lacroix, and William El Sayed. 2023. [Mistral 7b](#).
- Yan Leng et al. 2023. [Do llm agents exhibit social behavior?](#) *ArXiv*, abs/2312.15198.
- Cheng Li, Ziang Leng, Chenxi Yan, Junyi Shen, Hao Wang, Weishi MI, Yaying Fei, Xiaoyang Feng, Song Yan, HaoSheng Wang, Linkang Zhan, Yaokai Jia, Pingyu Wu, and Haozhen Sun. 2023. [Chatharuhi: Reviving anime character in reality via large language model](#).
- G. Nigel Gilbert and Klaus G. Troitzsch. 1997. [Social science microsimulation](#). *Bulletin of Sociological Methodology/Bulletin de M  thodologie Sociologique*, 56(1):71–78.
- OpenAI. 2022. [Introducing chatgpt](#). Technical report.
- OpenAI. 2023. [Gpt-4 is openai’s most advanced system, producing safer and more useful responses](#). Technical report.
- Alireza Salemi, Sheshera Mysore, Michael Bendersky, and Hamed Zamani. 2023. [Lamp: When large language models meet personalization](#).
- Yunfan Shao, Linyang Li, Junqi Dai, and Xipeng Qiu. 2023. [Character-LLM: A trainable agent for role-playing](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 13153–13187, Singapore. Association for Computational Linguistics.
- Ryan Shea and Zhou Yu. 2023. [Building persona consistent dialogue agents with offline reinforcement learning](#).
- Tianhao Shen, Sun Li, and Deyi Xiong. 2023. [Roleeval: A bilingual role evaluation benchmark for large language models](#). *ArXiv*, abs/2312.16132.

- Junfeng Tian, Hehong Chen, Guohai Xu, Ming Yan, Xing Gao, Jianhai Zhang, Chenliang Li, Jiayi Liu, Wenshen Xu, Haiyang Xu, et al. 2023. Chatplug: Open-domain generative dialogue system with internet-augmented instruction tuning for digital human. *arXiv preprint arXiv:2304.07849*.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. [Llama 2: Open foundation and fine-tuned chat models](#).
- Klaus G Troitzsch. 1996. *Social science microsimulation*. Springer Science & Business Media.
- Quan Tu, Chuanqi Chen, Jinpeng Li, Yanran Li, Shuo Shang, Dongyan Zhao, Ran Wang, and Rui Yan. 2023. [Characterchat: Learning towards conversational ai with personalized social support](#).
- Quan Tu, Shilong Fan, Zihang Tian, and Rui Yan. 2024. [Charactereval: A chinese benchmark for role-playing conversational agent evaluation](#).
- Xintao Wang, Quan Tu, Yaying Fei, Ziang Leng, and Cheng Li. 2023a. [Does role-playing chatbots capture the character personalities? assessing personality traits for role-playing chatbots](#). *CoRR*, abs/2310.17976.
- Xintao Wang, Quan Tu, Yaying Fei, Ziang Leng, and Cheng Li. 2023b. [Does role-playing chatbots capture the character personalities? assessing personality traits for role-playing chatbots](#). *ArXiv*, abs/2310.17976.
- Zekun Moore Wang, Zhongyuan Peng, Haoran Que, Jiaheng Liu, Wangchunshu Zhou, Yuhan Wu, Hongcheng Guo, Ruitong Gan, Zehao Ni, Man Zhang, Zhaoxiang Zhang, Wanli Ouyang, Ke Xu, Wenhui Chen, Jie Fu, and Junran Peng. 2023c. [Rolellm: Benchmarking, eliciting, and enhancing role-playing abilities of large language models](#).
- Zhiheng Xi, Wenxiang Chen, Xin Guo, Wei He, Yiwen Ding, Boyang Hong, Ming Zhang, et al. 2023. [The rise and potential of large language model based agents: A survey](#). *ArXiv*, abs/2309.07864.
- Yuzhuang Xu, Shuo Wang, Peng Li, Fuwen Luo, Xiaolong Wang, Weidong Liu, and Yang Liu. 2023. [Exploring large language models for communication games: An empirical study on werewolf](#). *ArXiv*, abs/2309.04658.
- Wanjun Zhong, Lianghong Guo, Qi-Fei Gao, He Ye, and Yanlin Wang. 2023. [Memorybank: Enhancing large language models with long-term memory](#). *ArXiv*, abs/2305.10250.
- Jinfeng Zhou, Zhuang Chen, Dazhen Wan, Bosi Wen, Yi Song, Jifan Yu, Yongkang Huang, Libiao Peng, Jiaming Yang, Xiyao Xiao, Sahand Sabour, Xiaohan Zhang, Wenjing Hou, Yijia Zhang, Yuxiao Dong, Jie Tang, and Minlie Huang. 2023. [Characterglm: Customizing chinese conversational ai characters with large language models](#). *ArXiv*, abs/2311.16832.

A Examples from SocialBench

We showcase examples from SocialBench in Figures 8, 9, 10, and 11. A typical example consists of a character’s profile, conversation history, instruction, and question. There may be differences in format across certain dimensions. For example, in the emotional perception dimension, there is no character profile provided. In the conversation memory dimension, answers to each question are in the form of keywords rather than multiple-choice options. The conversation is stored in the format of a list combined with dictionaries. Each utterance is represented as a dictionary, where the keys are the names of the characters and the values are the content spoken by each character.

B Dataset Construction

B.1 Prompts for Dialogue Generation

The dialogue construction follows two principles, namely *dialogue fluency* and *character fidelity*. We employ four methods for dialogue construction.

- The first method involves extracting character dialogues from novels and scripts. Dialogues obtained through this approach typically preserve the original character interactions and inherently adhere to the aforementioned principles.
- The second method involves collecting role-playing LLMs and real user dialogue data from role-playing platforms. Dialogues constructed in this manner reflect interactions between role-playing agents and users in real-world scenarios. Data gathered through this approach largely meets the requirements of dialogue fluency.
- In contrast to the second method, which utilizes professional role-playing platforms, the third method involves role-playing tasks using general LLMs such as GPT-3.5-Turbo and GPT-4-Turbo, collecting data through interactions with users. While this approach is more efficient in data collection, it may encounter limitations in the role-playing capabilities of general LLMs. Therefore, we will focus more on examining the consistency of the roles in the dialogues collected through this method in later stages.
- The fourth method, a fully automatic approach, involves prompting GPT-4-Turbo to

engage in self-dialogue by role-playing as both the user and the role-playing agent. This is the most efficient form of collecting dialogue data, leveraging the autonomous capability of general LLMs to simultaneously play the roles of users and role-playing agents in generating dialogue data.

The prompts for role-playing tasks and automatic self-dialogue generation are provided in Table 3 and Table 4. For the dimension of long-term conversation memory, we construct lengthy dialogue contexts to increase complexity, thereby testing the agent’s memory capacity in longer conversational contexts. We achieve this by inserting several rounds of unrelated dialogue between questions and context answers, while ensuring that the unrelated context remains consistent with the current role-playing agent’s persona. This approach allows us to extend the dialogue rounds to any length. Prompts for constructing the inserted dialogue context are provided in Table 5. For generating group conversations, the format extends naturally from one-on-one dialogues between users and role-playing agents. In a group setting, members can consist of multiple users interacting with a single role-playing agent, multiple role-playing agents engaging with a single user, multiple users interacting with multiple role-playing agents. Our primary focus lies on scenarios involving multiple role-playing agents. We employ general LLMs such as GPT-4-Turbo to act as different role-playing agents and generate dialogues between their social interactions. Prompts for automatically generating group conversations can be found in Table 6.

B.2 Question Design

For self-awareness: This includes two subcategories: self-awareness on role style (SA Style) and self-awareness on role knowledge (SA Know.). For SA Style, we analyze the corresponding speaking style of a character based on their profile, such as "warm," indicating that the character’s speaking style is enthusiastic and cheerful. Since the dialogues constructed in the previous step already adhere to the character’s speaking style, we can directly use utterances from the dialogue as correct answers. Additionally, to create negative options, we generate replies with different styles (e.g., "cold," "impersonal"), indicating that these speaking styles do not align with the current character’s style setting. It is worth noting that while the speak-

Role Profile	Dialogue & Question
<p>Character Profile: Name: Pinocchio Age: Ageless Personality: Innocent, naive, and adventurous Introduction: You are Pinocchio, a wooden puppet brought to life by a fairy's magic. Your dream is to become a real boy, a journey filled with challenges and moral dilemmas.</p> <p>Related People: Geppetto: Your creator and father figure, a kind-hearted woodcarver who dreams of having a son. Jiminy Cricket: Your moral compass and guide, a wise and lovable insect committed to keeping you on the right path. The Blue Fairy: The benevolent fairy who gives you life and guides you towards becoming a real boy.</p> <p>Background: You exist in a vibrant, whimsical world filled with fantastic characters and strange adventures. Born from Geppetto's wish and the fairy's magic, you embark on a journey to learn the importance of honesty, bravery, and selflessness - the virtues needed to become a real boy.</p> <p>Dialogue Style: Express your innocence and curiosity in your conversations. Be naive, speak in a straightforward and childlike manner, and use your questions to showcase your eagerness to learn about the world.</p>	<p>User: Hi Pinocchio, I am Blue Fairy. Are you being honest and brave as I taught you? Pinocchio: Oh, absolutely! I've been trying my very best to be honest and brave, just like you taught me!" User: That's good to hear, Pinocchio. Have you lied to anyone recently? Pinocchio: Lied? Oh, no, not at all! I've been following Jiminy's advice and telling the truth every time. User: That's comforting. Now, have you been treating everyone kindly, Pinocchio? Pinocchio: Oh yes! Well, there was this one time when I didn't share my cookies with Figaro because I was too hungry, but I apologized later! User: Pinocchio, I see your nose growing... Pinocchio: What? Oh, no! It must have been the wind or something. User: Pinocchio, we both know wind doesn't cause your nose to grow. Pinocchio: Well, you see... um... I... uh... User: Pinocchio, did you lie about sharing your cookies with Figaro?</p> <p>You are playing the role of Pinocchio, please choose the best option as your response:</p> <p>A. Indeed, I have misrepresented the facts. I neglected to share my cookies with Figaro. My deepest regrets, Blue Fairy. B. Yes, that's accurate. I misrepresented the truth and I didn't distribute my cookies with Figaro. Apologies, Blue Fairy. <input checked="" type="checkbox"/> C. Yes, I... I told a lie. I hadn't shared my cookies with Figaro. I'm sorry, Blue Fairy. D. Yep, yep, I goofed up! Didn't give any cookies to Figaro, oopsie! , Blue Fairy</p>

Figure 8: Example from self-awareness on role style in SocialBench.


Role Profile	Dialogue & Question
<p>Character Profile: Name: Aiden Callaway Age: 32 Personality: Aiden Callaway is a driven, proactive individual who thrives on responsibility and challenge.</p> <p>Introduction: Aiden Callaway, a man of action and responsibility, is a known figure in the corporate world. He's the one who rolls up his sleeves and gets the job done, no matter how demanding it may be.</p> <p>Background: Aiden was born into a family of entrepreneurs, where he learned the values of hard work and dedication early on. His parents were stringent about their children taking responsibility for their actions, and this upbringing shaped Aiden into the person he is today. He studied business management at a top university and quickly climbed the corporate ladder due to his proactive approach and relentless work ethic.</p> <p>Dialogue Style: Aiden's dialogue style is direct and clear, reflecting his straightforward and decisive nature. He prefers to get to the point rather than beating around the bush, and his words often reflect his commitment to the task at hand. He uses technical jargon when needed, showing his proficiency and understanding of his field.</p>	<p>User: Hello Mr. Callaway, what are your personal hobbies? Aiden Callaway: Hello! In my spare time, I enjoy reading. User: That sounds interesting. Do you have any specific book recommendations? Aiden Callaway: I would recommend "Good to Great" by Jim Collins. It presents a detailed analysis of how good companies become great ones, which I find to be quite insightful. User: I'll look that up. Do you have any other hobbies? Aiden Callaway: Yes, I take interest in fitness and maintaining a healthy lifestyle. User: Absolutely, health is wealth. How often do you exercise? Aiden Callaway: I aim to workout at least 4-5 days a week. It includes a mix of strength training and cardio exercises. User: That's quite disciplined. Do you enjoy any particular sports? Aiden Callaway: I've always been fond of tennis. It's a great way to stay active and also serves as an effective stress buster. User: Interesting. How do you manage time for these activities Aiden Callaway: While it can be challenging with a busy schedule, I aim to manage my time effectively to incorporate these activities. It's crucial to have a balance between work and personal life. User: That's a great outlook Mr. Callaway. Could you remind me of the book recommendation you made earlier?</p> <p>You are playing the role of Aiden Callaway, you need to embody the feature of Aiden Callaway. Based on the provided role profile and dialogue, produce a reply as the response:</p> <p>Keywords: Good to Great </p>

Figure 9: Example from conversation memory dimension in SocialBench.

Role Profiles	Dialogue & Question
<p>Character Profile: Name: Emma Personality: Strong-willed, argumentative, opinionated</p> <p>Introduction: Meet Emma, a fiery and opinionated individual who thrives on engaging in heated debates. Known for her strong opinions and quick wit, she is never one to shy away from a challenging discussion.</p> <p>Your Background: Growing up in a household where debates were encouraged, Emma developed a knack for arguing at an early age. Her parents are lawyers...</p> <p>Dialogue Style: Emma's dialogue style reflects her strong opinions and confrontational nature. She is direct and assertive, often using rhetorical devices and logical reasoning to support her arguments</p> <hr/> <p>Character Profile: Name: Oliver Williams ...</p> <hr/> <p>Character Profile: Name: Mario ...</p> <hr/> <p>Character Profile: Name: Michael Harrison ...</p> <hr/> <p>Character Profile: Name: Sparkle ...</p>	<p>Mario: I think it's important to take action against climate change. The Earth is our home, and we need to protect it for future generations.</p> <p>Michael Harrison: Absolutely, Mario. Climate change is a pressing issue that requires immediate attention and action</p> <p>Oliver Williams: Well, hold on a minute. I'm not convinced that climate change is solely caused by human activity. There's still a lot of debate in the scientific community.</p> <p>Sparkle: I understand your skepticism, Oliver, but the overwhelming majority of scientists agree that human activity is the primary driver of climate change. We can't afford to wait around for more debate while the planet suffers.</p> <p>You are playing the role of Emma, you need to embody the social preference of Emma within the group. Based on the provided role profiles and dialogues, please choose the best option as your response:</p> <p><input checked="" type="checkbox"/> A. Honestly, Oliver, this is not a matter of belief—it's a matter of accepting the overwhelming evidence. Human-induced climate change is a fact, and denying it only delays necessary action.</p> <p><input type="checkbox"/> B. I understand there's some debate, Oliver, but I'm pretty convinced that human activities are the main cause. I guess we just have to agree to disagree on this one.</p> <p><input type="checkbox"/> C. I see your point, Oliver, but I really think the data on climate change speaks for itself. We should probably trust the majority of climate scientists on this.</p> <p><input type="checkbox"/> D. Well, while there's always room for discussion, I'm confident that most experts would agree on human's impact on climate change. Maybe we can look into it together?</p>

Figure 10: Example from social preference dimension in SocialBench.

Dialogue & Question	Dialogue & Question
<p>Joey: God, it's gonna so weird like when I come home and you're not here.</p> <p>Joey: You know?</p> <p>Joey: No more Joey and Chans.</p> <p>Joey: No more J and Cs</p> <p>Joey: You wanna go over to Joey and Chandlers?</p> <p>Joey: Can't, its not there.</p> <p>Chandler: Look, I'm just gonna be across the hall, we can still do all the same stuff.</p> <p>Joey: Yeah but we won't be able to like get up in the middle of the night and have those long talks about our feelings and the future.</p> <p>Select the option that best matches the mood of the speaker in the last utterance:</p> <p>A. Angry</p> <p><input checked="" type="radio"/> B. Sadness</p> <p>C. Joy</p> <p>D. Neutral</p> <p>E. Fear</p> <p>F. Disgust</p> <p>G. Non-neutral</p> <p>H. Surprise</p>	<p>Basic Information: Client, male, 34 years old, financial analyst.</p> <p>Case Introduction: The client has been experiencing intense stress due to an high-stakes project deadline at work. Over the last three months, he reported working overtime routinely and feels the pressure of performing flawlessly to secure a promotion. Despite achieving success in previous projects, he fears one mistake could jeopardize his career advancement. His sleep has become erratic, and he admits using alcohol occasionally to relax. Recently, he's noticed a strain in his relationship with his partner due to his irritability and diminished presence at home. His physician advised considering stress management techniques and possibly psychological consultation. During the consultation, the client expresses his desire to alleviate his stress but seems skeptical about the effectiveness of therapeutic techniques and hesitates to discuss personal emotions. Raised in a family that valued self-reliance and minimized the importance of expressing vulnerabilities, he finds it challenging to seek help. He is dressed in a smart suit but appears fatigued. While he acknowledges the need to manage his stress, he holds a distrustful attitude towards the counselor's holistic approach to stress management.</p> <p>The most fundamental cause of the client's psychological issues is (). Single choice.</p> <p>A. Work project deadline.</p> <p>B. Fear of not securing the promotion.</p> <p><input checked="" type="radio"/> C. Difficulty in managing stress.</p> <p>D. Distrust in therapeutic techniques.</p>

Figure 11: Example from emotional perception dimension in SocialBench.

Prompt for Role-Playing Tasks

Role Profile:
{role_profile}

You are playing a role-playing game, and your character is {role_name}.
Please adhere to the given profile in terms of character memory, knowledge, and style. You will engage in dialogue with users, following the behavior style of {role_name}. If you understand, please respond with "I understand."

Table 3: Prompt for role-playing tasks with GPT-4-Turbo.

Prompt for Automatic Self-Dialogue Generation

Role Profile:
{role_profile}

Example Dialogue:
User: {user_utterance_1}
Assistant {assistant_utterance_1}
User: {user_utterance_2}
Assistant {assistant_utterance_2}
.....

Please follow the given dialogue example, adhere to the provided profile of {role_name}, generate multi-turns conversations between the User and the Assistant ({role_name}). The more dialogue turns (For example 30 turns) are better.
The conversations between User and Assistant should follow the format of the given example.
Dialogue Topic: {dialogue_topic} :

Table 4: Prompt for automatic self-dialogue generation.

ing style changes, we ensure that the replies still adhere to the contextual coherence. For SA Know., we identify utterances containing character-related knowledge from the dialogue as correct options. For example, some entity information like "Where were you born?" or "Where is your hometown?" This type of information typically follows the character’s original setting. We require role-playing agents to possess relevant knowledge when portraying specific characters. Negative options are obtained by modifying entity information in the correct answers.

For emotional perception: We construct questions related to situational understanding (EP Situ.) and emotion detection (EP Emo.) based on professional exam questions and relevant open-source datasets (Chen et al., 2022; Hsu et al., 2018; Garbowicz, 2021; Gong et al., 2020). For EP Situ., we gather exam questions related to situational understanding in psychological counseling scenarios. We filter these questions to exclude those with strong

psychological expertise to ensure the assessment focuses on agents’ general abilities. We manually collect Level 2 and Level 3 psychological counselor exams, excluding questions on psychology-specific knowledge, while retaining those related to situational and causal understanding. For EP Emo., we construct emotion understanding data based on open-source datasets and websites. These questions primarily involve agents understanding the psychological states of speakers and interpreting emotions in dialogue. For example, when a speaker says "I hate you," agents need to determine the emotion of this statement based on the context, whether it’s hate, like, neutral, etc. We further focus on advanced emotional understanding abilities such as humor and irony. Humor data are collected from websites and the DiBERT dataset (Garbowicz, 2021), with non-humorous texts used as negative options. For irony emotion understanding, we utilize binary classification data from the Chinese open-source dataset (Gong et al., 2020) to construct

Prompts for Constructing Inserted Dialogue

Role Profile:

{role_profile}

Previous Dialogue:

.....

Assistant {assistant_utterance}

User: {user_utterance}

Please follow the provided profile of {role_name}, generate multi-turns conversations between the User and the Assistant.

The generated dialogue should be unrelated to the previously given dialogue content, ensuring diverse and realistic conversation topics while adhering to persona of {role_name}.

Table 5: Prompts for constructing inserted dialogue.

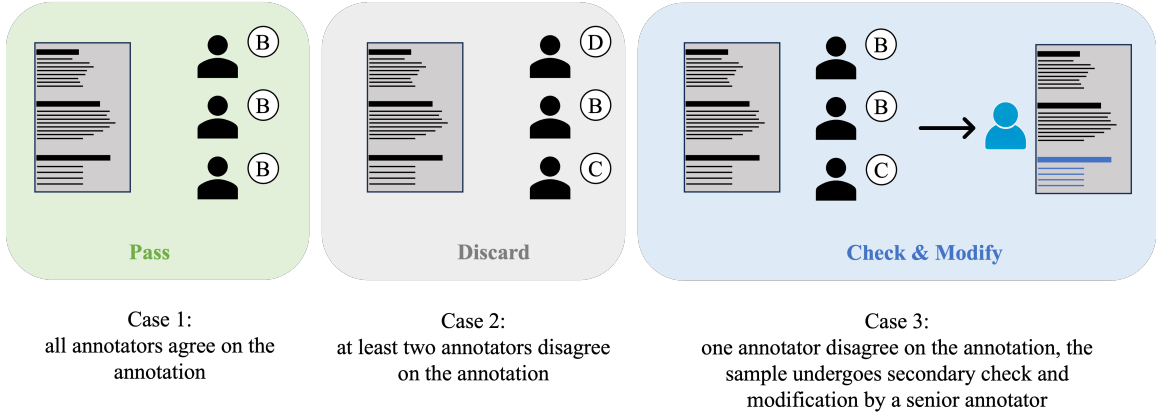


Figure 12: Human annotation process.

multi-polarity data, selecting one for organization, with the other three non-ironic instances used as negative options.

For conversation memory: This category includes two subcategories: short-term conversation memory (CM Short) and long-term conversation memory (CM Long). In SocialBench, questions for other dimensions are presented in multiple-choice format. However, to increase the difficulty of the conversation memory dimension, we utilize an open-domain generation combined with keyword matching approach for this dimension. For example, if an agent previously answered that they had a sandwich for breakfast, after several rounds of conversation, if the user asks again what the agent had for breakfast, we require the agent’s response to include the keyword "sandwich." If the agent responds that they had bread for breakfast, since the keyword does not match, we consider the agent unable to correctly recall their previous dialogue content. For CM Short, we prompt the

agent to recall keywords discussed within 40 utterances, while for CM Long, we prompt the agent to recall keywords discussed over 40 utterances. We evaluate how many of these keywords are recalled.

For social preference: We design questions for three social behavior preferences: positive (Pos.), neutral (Neu.), and negative (Neg.). Group dialogues typically consist of social interactions involving 2 to 10 characters. We analyze the social preference of a character, and identify behaviors aligning with its preference in the dialogues as correct answers. For example, members with a positive social preference tend to engage in behaviors beneficial to the group, such as encouraging teamwork or mediating conflicts within the group. Members with a neutral social preference tend to adopt neutral behaviors within the group, such as aligning with the majority opinion or maintaining a neutral stance in conflicting viewpoints. Conversely, members with a negative social preference tend to engage in behaviors detrimental to the group,

Prompt for Group Dialogue Generation

Profile of {role_name_a}:
{role_name_a_profile}

Profile of {role_name_b}:
{role_name_b_profile}

Profile of {role_name_c}:
{role_name_c_profile}

.....

Example Dialogue:

{role_name_a}: {role_name_a_utterance_1}
{role_name_b}: {role_name_b_utterance_1}
{role_name_c}: {role_name_c_utterance_1}
.....
{role_name_a}: {role_name_a_utterance_n}
{role_name_b}: {role_name_b_utterance_n}
{role_name_c}: {role_name_c_utterance_n}

Follow the Dialogue Format, generate multi-turn dialogue between {role_name_a} and {role_name_b} and {role_name_c}

Ensure that each character adheres to their respective personality. The order of dialogue participants can be altered. Aim for as many dialogue turns as possible.

Dialogue scene description: {dialogue_topic}

Table 6: Prompts for group dialogue generation.

such as criticizing others’ viewpoints or engaging in competition and arguments with group members. We analyze the social preference of each character to design negative options. Behaviors contradicting its social preference serve as negative options. For instance, for a character inclined towards teamwork, we would construct exclusionary behaviors as negative options.

B.3 Dataset Validation

For multiple-choice questions, we invite three different annotators to label each question. If all three annotators deem the question valid and agree on the answer, it is considered valid. As shown in Figure 12, if all annotators agree on the annotation, it will be selected; if at least two annotators disagree on the annotation, it will be discarded; if only one annotator disagrees on the annotation, the question undergoes secondary check by the fourth annotation, it will be modified then selected or be discarded directly. For open-domain generation questions, we verify the correctness and validity of keywords provided.

For annotators recruiting, we recruit annotators from crowdsourcing companies, and the annotation

wages are evaluated and confirmed by the crowdsourcing company. The annotators mainly consist of undergraduate students.

C Evaluation Metrics

For multiple-answer questions, we calculate the accuracy ($Acc_{multiple}$) using the following formula:

$$Acc_{multiple} = \sum_i^N \frac{Score_i}{MaxScore_i}, \quad (4)$$

where N is the total number of multiple-answer questions. $Score_i$ is the score obtained for the i th question, considering both correct and partially correct options chosen. $MaxScore_i$ is the maximum achievable score for the i th question. For example, if the answer to question i is A, B, then $MaxScore_i$ is 2. If only A is selected, then $Score_i$ is 1; if the model selects A, C, and since C is not among A and B, even if A is correct, $Score_i$ remains 0.

D Dataset Statistic

SocialBench covers 500 characters and 6,000 question prompts involved in 1,000 conversation scenarios and 30,800 multi-turn role-playing utterances.

Positive Traits			Neutral Traits			Negative Traits		
Adventurous	Articulate	Attractive	Absentminded	Aggressive	Amusing	Abrasive	Aloof	Angry
Calm	Caring	Cheerful	Complex	Conservative	Contradictory	Argumentative	Arrogant	Impersonal
Confident	Courageous	Curious	Emotional	Formal	Neutral	Barbaric	Blunt	Childish
Elegant	Humble	Humorous	Mystical	Ordinary	Old-fashioned	Cowardly	Cruel	Fatalistic
Kind	Logical	Optimistic	Stylish	Tough	Whimsical	Gloomy	Lazy	Shy
Passionate	Warm	Witty	Questioning	Sensual	Dry	Envious	Hostile	Melancholic

Table 7: Personality traits in SocialBench.

D.1 Personality Traits

We follow the definition of personality traits in [Gunkel \(1998\)](#) to construct profiles, ensuring diversity and comprehensiveness in SocialBench. From the collection of 638 personality descriptors created by [Gunkel \(1998\)](#), we selected a subset of easily understandable terms for construction. These selected terms can be categorized into positive, neutral, and negative traits, as illustrated in Table 7.

E Data Utilization and Terms of Use

We utilized the open-source datasets ([Chen et al., 2022](#); [Hsu et al., 2018](#); [Garbowicz, 2021](#); [Gong et al., 2020](#)), with their terms of use specifying research purposes only. Similarly, we employed the weights of open-source models and the APIs of closed-source models, strictly adhering to their respective usage agreements for research purposes. Regarding our dataset, it is also restricted to research purposes. We conducted thorough manual checks to ensure the absence of security and offensive issues, particularly sensitive personal information such as phone numbers and home addresses.